

Comparison of selected model order reduction methods

Studienarbeit (T3_3101)

im Rahmen der Prüfung zum
Bachelor of Science (B.Sc.)

des Studienganges Angewandte Informatik
an der Dualen Hochschule Baden-Württemberg Karlsruhe

von

Florian Braun

Januar 2018

-Sperrvermerk-

Abgabedatum:	22. Mai 2023
Bearbeitungszeitraum:	14.10.2022 - 22.05.2023
Matrikelnummer, Kurs:	7433149, TINF20B1
Gutachter der Dualen Hochschule:	Lutz Gröll

Eidesstattliche Erklärung

Ich versichere hiermit, dass ich meine Studienarbeit (T3_3101) mit dem Thema:

Comparison of selected model order reduction methods

gemäß § 5 der "Studien- und Prüfungsordnung DHBW Technik" vom 29. September 2017 selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Die Arbeit wurde bisher keiner anderen Prüfungsbehörde vorgelegt und auch nicht veröffentlicht.

Ich versichere zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

Karlsruhe, den March 28, 2023

Nachname, Vorname

Abstract

- Deutsch -

Dies ist der Beginn des Abstracts. Für die finale Bachelorarbeit musst du ein Abstract in deinem Dokument mit einbauen. So, schreibe es am besten jetzt in Deutsch und Englisch. Das Abstract ist eine kurze Zusammenfassung mit ca. 200 bis 250 Wörtern.

Versuche in das Abstract folgende Punkte aufzunehmen: Fragestellung der Arbeit, methodische Vorgehensweise oder die Hauptergebnisse deiner Arbeit.

Contents

Abkürzungsverzeichnis	IV
List of Figures	V
List of Tables	VI
source-codeverzeichnis	VII
1. Fundamentals	1
1.1. Heat Equation	1
1.2. Finite Element Method	3
1.3. Singular Value Decomposition	6
1.4. Fundamentals of Control Theory	10
2. Model Order Reduction	13
2.1. Introduction	13
2.2. Proper Orthogonal Decomposition	14
2.3. Balanced Truncation	17
3. Implementation	20
3.1. Class main	21
3.2. Class containerFEM	21
Literaturverzeichnis	VIII
A. Appendix	I
A.1. Deriving matrices for FEM using piecewise linear functions	I
A.2. Proof that matrix M is invertible	III
A.3. Equivalence of picewise linear polynomials and linear interpolation	V

List of abbreviations

DE	Differential Equation
DFT	Discrete Fourier Transform
FEM	Finite Element Method
FT	Fourier Transform
FFT	Fast Fourier Transform
IDFT	Inverse Discrete Fourier Transform
IFFT	Inverse Fast Fourier Transform
LS	Least Squares
ODE	Ordinary Differential Equation
PDE	Partial Differential Equation

List of Figures

1.1. Variance and commulative variance captured by each column vectors of U and V	9
3.1. Class diagramm of MOR and FEM implementation	20
3.2. Flow chart of main class	21
3.3. Flow chart for FEM class	21

List of Tables

source-codeverzeichnis

1. Fundamentals

1.1. Heat Equation

The conduction of heat within a medium can be described using the following partial differential equation (PDE):

$$\frac{\partial u}{\partial t} = \alpha \nabla^2 u \quad (1.1)$$

With u being a function of space and time and α being a positive constant. For this paper u will be defined in terms of one spacial dimension:

$$u := u(x, t) \quad (1.2)$$

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2} \quad (1.3)$$

$$x \in \chi \subset \mathbb{R} \quad t \in \tau \subset \mathbb{R} \quad (1.4)$$

$$x_0 \leq x \leq x_n \quad t_0 \leq t \leq t_n \quad (1.5)$$

[1]

In order to not only model the conduction of heat within a medium but also a heating process, a new function $h : \chi \times \tau \rightarrow \mathbb{R}$ is introduced:

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2} + h(x, t) \quad (1.6)$$

For this paper it is assumed that the initial condition is known:

$$f : \chi \rightarrow \mathbb{R} \quad (1.7)$$

$$u(x, t_0) = f(x) \quad (1.8)$$

1.1.1. Solving the Heat Equation using Fourier transform

Applying the Fourier transform (FT) w.r.t x to 1.6 yields the inhomogeneous ordinary differential equation (ODE):

$$\hat{u} = \mathfrak{F}(u) \quad \hat{h} = \mathfrak{F}(h) \quad (1.9)$$

$$\frac{d}{dt} \hat{u} = -\alpha \omega^2 \hat{u} + \hat{h} \quad (1.10)$$

A solution to 1.10 is given by:

$$\hat{u} = \hat{u}_0 + \hat{u}_p \quad (1.11)$$

Where \hat{u}_0 is the homogeneous solution and \hat{u}_p is the particular integral. In order to solve this ODE for the particular integral \hat{h} has to be known. [2] The choice of h is, except to some restrictions, arbitrary. Therefore an approximate solution to 1.10 \hat{u}_a is obtained by the forward euler scheme:

$$\frac{d}{dt}\hat{u} \approx \frac{\Delta\hat{u}}{\Delta t} \quad (1.12)$$

$$\hat{u}_{t+1} = \hat{u}_t + \Delta t(-\alpha\omega^2\hat{u} + \hat{h}) \quad (1.13)$$

$$\hat{u}_a = [\hat{u}_{t_0}, \dots, \hat{u}_{t_n}] \quad (1.14)$$

In order to apply the euler scheme successfully an initial condition \hat{u}_0 has to be known. This initial condition is obtained by applying the discrete Fourier transform (DFT) to an initial temperature distribution along x :

$$\hat{u}_0 = \mathfrak{F}(f(x)) \quad (1.15)$$

[3]

The forward euler scheme is used here because it is fairly easy to implement. By applying the inverse discrete Fourier transform (IDFT) to \hat{u}_a an approximate solution to 1.6 can be obtained. To decrease computing time, the fast Fourier transform (FFT) and inverse fast Fourier transform (IFFT) is used instead of the DFT and IDFT.

1.2. Finite Element Method

The finite element method (FEM) is a method to approximate solutions for differential equations (DE) within a certain domain Ω . This is done by discretizing the spacial domain. Assume that a DE is given by:

$$m, n \in \mathbb{N} \quad \zeta \in \Omega \subset \mathbb{R} \quad m \geq 1 \quad (1.16)$$

$$\frac{\partial^m y}{\partial \zeta^m} - g(y) = r(\zeta, t) \quad (1.17)$$

It is assumed that g is a linear function that can also contain partial derivatives of y w.r.t. time, y takes the value 0 at the boundary Γ and $y(\zeta, 0) = f(\zeta)$. An approximate solution to y is given by μ , which is expressed as a sum of basis functions contained in the set ϕ :

$$\mu(\zeta, t) = \sum_{j=1}^N c_j(t) \phi_j(\zeta) \quad (1.18)$$

The residual is defined as:

$$\mathbf{r} = \frac{\partial^m \mu}{\partial \zeta^m} - g(\mu) - r(\zeta, t) \quad (1.19)$$

Furthermore the residual is required to be orthogonal to all basis functions:

$$\langle \mathbf{r}, \phi_k \rangle = 0 \quad \forall \phi_k \in \phi \quad (1.20)$$

Since the functions in ϕ are known, it is only required to find the coefficients $c_j(t)$ in 1.18. To find those coefficients 1.20 needs to be expressed as follows:

$$\int_{\Omega} \frac{\partial^m \mu}{\partial \zeta^m} \phi_k d\zeta - \int_{\Omega} g(\mu) \phi_k d\zeta = \int_{\Omega} r(\zeta, t) \phi_k d\zeta \quad \forall \phi_k \in \phi \quad (1.21)$$

If μ is substituted with 1.18 the following is obtained:

$$\sum_{j=1}^N \left(\left(\int_{\Omega} \frac{\partial^m \phi_j}{\partial \zeta^m} \phi_k d\zeta \right) c_j(t) - g \left(\left(\int_{\Omega} \phi_k \phi_j d\zeta \right) c_j(t) \right) \right) = \int_{\Omega} r(\zeta, t) \phi_k d\zeta \quad \forall \phi_k \in \phi \quad (1.22)$$

It is also necessary to apply divergence theorem to the first integral term taking into account that y at Γ is 0. Since ζ is one dimensional, the divergence theorem becomes integration by parts:

$$\int_{\Omega} \frac{\partial^m \phi_j}{\partial \zeta^m} \phi_k d\zeta = - \int_{\Omega} \frac{\partial^{m-1} \phi_j}{\partial \zeta^{m-1}} \frac{\partial \phi_k}{\partial \zeta} d\zeta \quad \forall \phi_k \in \phi \quad (1.23)$$

Combining 1.22 and 1.23 yields:

$$-\sum_{j=1}^N \left(\int_{\Omega} \frac{\partial^{m-1} \phi_j}{\partial \zeta^{m-1}} \frac{\partial \phi_k}{\partial \zeta} d\zeta \right) c_j(t) + g \left(\int_{\Omega} \phi_k \phi_j d\zeta \right) c_j(t) = \int_{\Omega} r(\zeta, t) \phi_k d\zeta \quad \forall \phi_k \in \phi \quad (1.24)$$

This formulation leads to a system of ODEs or a system of linear equations that can be solved either analytically or numerically.

1.2.1. Solving the Heat Equation using FEM

This formulation of FEM can be applied to 1.6:

$$\Omega = \chi \quad \Gamma = \{x_0, x_n\} \quad (1.25)$$

$$y(\zeta, t) = -u(x, t) \quad g(u) = -\frac{1}{\alpha} \frac{\partial u}{\partial t} \quad (1.26)$$

$$m = 2 \quad r(\zeta, t) = \frac{1}{\alpha} h(x, t) \quad (1.27)$$

$$u(x, 0) = f(x) \quad u(x_0, t) = 0 \quad u(x_n, t) = 0 \quad (1.28)$$

The set of basis functions is defined as a set of piecewise linear functions with constant step size Δx :

$$\phi_j(x) = \begin{cases} (x - x_{j-1})/\Delta x, & x_{j-1} \leq x < x_j \\ (x_{j+1} - x)/\Delta x, & x_j \leq x < x_{j+1} \\ 0, & \text{otherwise} \end{cases} \quad (1.29)$$

[4]

The stepsize Δx is defined by $\Delta x = \frac{x_n - x_0}{n-1}$. This results in the following system of ODEs:

$$\sum_{j=1}^N \left(\int_{\chi} \phi_j \phi_k dx \right) \frac{dc_j}{dt} = \alpha \sum_{j=1}^N \left(- \int_{\chi} \frac{d\phi_j}{dx} \frac{d\phi_j}{dx} dx \right) c_j(t) + \int_{\chi} h(x, t) \phi_k dx \quad \forall \phi_k \in \phi \quad (1.30)$$

Using matrix notation this becomes:

$$M^{N \times N}, K^{N \times N} \quad (1.31)$$

$$M \dot{c} = Kc + d \quad (1.32)$$

The matrices M and K can be easily computed (Appendix A.1):

$$m_{ij} = \begin{cases} \frac{2\Delta x}{3}, & k = j \\ \frac{\Delta x}{6}, & |k - j| = 1 \\ 0, & \text{otherwise} \end{cases} \quad k_{ij} = \begin{cases} \frac{-2\alpha}{\Delta x}, & k = j \\ \frac{\alpha}{\Delta x}, & |k - j| = 1 \\ 0, & \text{otherwise} \end{cases} \quad (1.33)$$

However it is necessary to approximate d for each point in time using numerical integration schemes. Furthermore to solve this system of ODEs numerically an initial condition c_0 has to be known. [3]

It can be obtained using a least squares (LS) approach:

$$\sum_{j=1}^N \langle \phi_j, \phi_k \rangle c_j(0) = \langle f, \phi_k \rangle \quad \forall \phi_k \in \phi \quad (1.34)$$

$$M c_0 = F \quad (1.35)$$

$$c_0 = M^{-1} F \quad (1.36)$$

[5]

Observe that by multiplying 1.32 with M^{-1} (A.2) yields a system of ODEs:

$$\dot{c} = M^{-1} K c + M^{-1} d \quad (1.37)$$

This system of ODEs can be solved using an euler scheme:

$$c_{t+1} = \Delta t M^{-1} (K c_t + d) + c_t \quad (1.38)$$

Force boundary conditons

Keep in mind that vector d is time dependent and has to be recomputed for each time step. To force the boudary condition $u(x, t) = 0 \quad x \in \Gamma$ the first and last entry of any c_t has to be zero:

$$c_t^1 = 0 \quad c_t^n = 0 \quad \forall t \quad (1.39)$$

Therefore the Matrix M^{-1} has to be adjusted:

$$0 = \Delta t m (K c + d) \quad (1.40)$$

$$\Rightarrow m = [0, 0, 0, \dots, 0] \quad (1.41)$$

Here m is the first row vector of M^{-1} or the last one respectively [3]. Using 1.18 and the computed coefficients c the function $u(x, t)$ can be approximated. However this is equivalent to linear interpolation between c_n and c_{n+1} (Appendix A.3).

1.3. Singular Value Decomposition

The Singular Value Decomposition (SVD) is a matrix factorization with guaranteed existence. It can be used to obtain low rank approximations of a matrix or pseudo inverses for ill posed linear system of equations. It is also related to FT by providing a data specific set of orthogonal bases instead of a generic set of sines and cosines. For this paper the SVD will be used for generating low rank approximations of matrices. [6]

1.3.1. Properties

A matrix $X \in \mathbb{C}^{n \times m}$ can be decomposed in the following way:

$$X = U \Sigma V^* \quad (1.42)$$

Here $U \in \mathbb{C}^{n \times n}$ and $V \in \mathbb{C}^{m \times m}$ are unitary matrices and $\Sigma \in \mathbb{R}^{n \times m}$ is a real valued ordered diagonal matrix. The columns of U provides a set of orthonormal basis vectors for the column space of X , V contains orthonormal basis vectors for the row space of X . The matrix Σ assigns a magnitude ('importance') to the product of U and V^* [7]. Since U and V are unitary they have the following property:

$$U^* U = U U^* = I \quad (1.43)$$

$$V^* V = V V^* = I \quad (1.44)$$

[8]

In case $n \geq m$ the so called economy SVD can be used to factorize the matrix X :

$$X = [\hat{U} \quad \hat{U}^\perp] \begin{bmatrix} \hat{\Sigma} \\ 0 \end{bmatrix} V^* = \hat{U} \hat{\Sigma} V^* \quad (1.45)$$

The economy SVD omits rows only containing zeros in Σ and the according columns of U . Therefore the dimensionality of \hat{U} and $\hat{\Sigma}$ is less or equal to the dimensionality of U and Σ . [6]

1.3.2. Hirachy of correlations

As already stated, the matrix Σ assigns a magnitude to UV^* . This magnitude can be seen as the square of the variance the bases in U and V capture. Assume that $X^* X$ and $X X^*$ denote a correlation matrix. [6] A correlation matrix is a matrix that stores correlation coefficients between multiple measurements. If X has the following properties $X^* X$ and $X X^*$ are correlation matrices:

1.) The column vectors of X have to be zero mean:

$$X = [x_1, \dots, x_m] \quad (1.46)$$

$$\frac{1}{n} \sum_{j=1}^n x_{ij} = \mu_i = 0 \quad \forall 1 \leq i \leq m \quad (1.47)$$

2.) The column vectors of X have to be normalized:

$$\left(\sum_{j=1}^n x_{ij}^2 \right)^{\frac{1}{2}} = \bar{x}_i = 1 \quad \forall 1 \leq i \leq m \quad (1.48)$$

The correlation coefficient between two column vectors of X is calculated as follows:

$$\text{corr}(x_i, x_{i'}) = \frac{(\sum_{j=1}^n x_{ij} - \bar{x}_i)(\sum_{j=1}^n x_{i'j} - \bar{x}_{i'})}{(\sum_{j=1}^n (x_{ij} - \bar{x}_i)^2)^{\frac{1}{2}} (\sum_{j=1}^n (x_{i'j} - \bar{x}_{i'})^2)^{\frac{1}{2}}} \quad (1.49)$$

[9]

Since all column vectors have zero mean and are normalized this becomes:

$$\text{corr}(x_i, x_{i'}) = \text{cov}(x_i, x_{i'}) = x_i^* x_{i'} \quad (1.50)$$

[10]

This resembles the entries of XX^* and X^*X . A vector e that maximizes the variance of the projections of x_i onto e with the restriction $\|e\| = 1$ are the eigenvectors of the according correlation matrix. The eigenvalue of e denoted as λ is equivalent to the variance of x_i projected onto e [11].

Since X can be deconstructed using the SVD, XX^* and X^*X are equal to:

$$XX^* = U \begin{bmatrix} \hat{\Sigma} \\ 0 \end{bmatrix} V^* V \begin{bmatrix} \hat{\Sigma} & 0 \end{bmatrix} U^* = U \begin{bmatrix} \hat{\Sigma}^2 & 0 \\ 0 & 0 \end{bmatrix} U^* \quad (1.51)$$

$$X^*X = V \begin{bmatrix} \hat{\Sigma} & 0 \end{bmatrix} U^* U \begin{bmatrix} \hat{\Sigma} \\ 0 \end{bmatrix} V^* = V \hat{\Sigma}^2 V^* \quad (1.52)$$

By multiplying U and V respectively on the right side 1.51 and 1.52 become:

$$XX^*U = U \begin{bmatrix} \hat{\Sigma}^2 & 0 \\ 0 & 0 \end{bmatrix} \quad (1.53)$$

$$X^*XV = V \hat{\Sigma}^2 \quad (1.54)$$

This shows that V contains the eigenvectors of the row-wise correlation matrix and U contains the eigenvectors of the column-wise correlation matrix. The matrix Σ contains the roots of the according eigenvalues and is thereby related to the variance. By ordering U , Σ and V by

the entries of Σ in a descending order the first row of U and V contain the most important basis vectors. [6]

1.3.3. Low-rank approximation

A usefull property of the SVD is that it can be used to find an hierachy of rank- r approximation for a given matrix X . An matrix \tilde{X} that approximates X is obtained by:

$$\tilde{X} = \arg \min \|X - \tilde{X}\|_F = \tilde{U} \tilde{\Sigma} \tilde{V}^* \quad (1.55)$$

$$s.t. rank(\tilde{X}) = r \quad (1.56)$$

Here \tilde{U} and \tilde{V} denote matrices obtained takeing the first r columns of U and V . The matrix $\tilde{\Sigma}$ is a $r \times r$ sub-block of Σ . This is alsow known as the Eckard-Young theorem. The variance captured by \tilde{X} can be calculated in the following way:

$$cumvar_r(\Sigma) = \frac{\sum_{i=1}^r \sigma_i}{trace(\Sigma)} \quad (1.57)$$

$$var(\Sigma) = \frac{diag(\Sigma)}{trace(\Sigma)} \quad (1.58)$$

[6]

1.3.4. Example low-rank approximation

As an example suppose there is a matrix $X \in \mathbb{Z}^{4 \times 4}$:

$$X = \begin{bmatrix} 3 & 1 & 5 & 5 \\ 4 & -4 & 5 & 0 \\ -4 & -2 & -4 & 3 \\ 5 & 1 & 5 & -4 \end{bmatrix} \quad (1.59)$$

By computing the SVD the matrices U , Σ and V are obtained. Now Σ can be used to calculate the cummulative variance for an rank- r approximation and the variance captured by each basis vector of U and V .

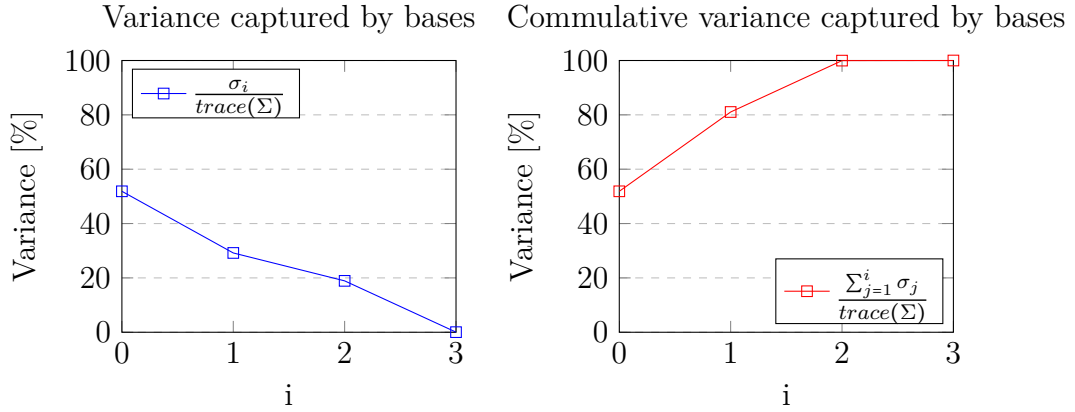


Figure 1.1.: Variance and commulative variance captured by each column vectors of U and V

On 1.3.4 the commulative variance and the variance of each basis vectors is plotted. Here X can be apprixomated using the first three leading basis vectors. This approximation captures already more than 99% of the variance. The resulting matrix \tilde{X}_3 looks as follows:

$$\tilde{X}_3 = \begin{bmatrix} 4.73 & -0.76 & 5.62 \\ 1.31 & -1.37 & 1.62 \\ -5.10 & -0.60 & 2.34 \end{bmatrix} \quad (1.60)$$

1.4. Fundamentals of Control Theory

1.4.1. Problems in Control Theory

There are mainly three different problems in control theory. To analyse a system it has to be identified first. A mathematical model has to be found that relates the input of that system to its output. After a model has been found, it is helpful to simulate the system to find out which input results in which output. Simulations have to fit the system and offer the advantage over real experiments of reducing cost and risks attached to them. The last of the three major problems is the control problem. Now that the underlying model is known and the system can be simulated, a controller for that system can be designed. Here the controller is a system that feeds some input to the system that is to be controlled to generate some desired output. [12] For this paper the second problem is the most important one, since the goal of this paper is to compare some methods to make simulations of a system more efficient.

1.4.2. LTI Systems

A linear and time invariant (LTI) system is a kind of system that fulfills the following conditions:

Homogeneity: For a system to be a LTI system, the condition that the output of a system $y(t)$ relates to the input $u(t)$ by a linear operator f :

$$y(t) = f(u(t)) \quad (1.61)$$

$$cy(t) = f(cu(t)) \quad (1.62)$$

This means that if the input u changes in magnitude by a constant factor of c the output y also changes by the factor of c .

Superposition: The second requirement is superposition:

$$y_1(t) + y_2(t) = f(u_1(t) + u_2(t)) \quad (1.63)$$

The sum of two outputs y_1 and y_2 has to equal the output of the system with the sum of the inputs u_1 and u_2 as input.

Time Invariance: A system that is time invariant has the property that if the input u is shifted in time by some constant τ the resulting output is also shifted in time by the same constant:

$$y(t - \tau) = f(u(t - \tau)) \quad (1.64)$$

LTI systems are important because they can be used to approximate non LTI systems over some region. This is useful since LTI systems are well understood. [13]

1.4.3. State Space representation

A LTI system can be expressed as a system of ODEs that relates the input u to its output y . The vector x denotes the states of a system.

$$\dot{x} = Ax + Bu \quad x(t_0) = x_0 \quad (1.65)$$

$$y = Cx + Du \quad (1.66)$$

The matrices $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{q \times n}$ and $D \in \mathbb{R}^{q \times m}$ are constant matrices. [14]

1.4.4. Transfer Function

To understand the transfer function the impulse response has to be defined first. The impulse response of a system is given by choosing the following function as input to the system:

$$g(t) = f(\delta(t)) \quad (1.67)$$

The function $\delta(t)$ is the delta function:

$$\delta(t) = \begin{cases} \infty, & t = 0 \\ 0 & \text{else} \end{cases} \quad (1.68)$$

It has the property that the integral of this function is 1:

$$\int_{t=-\infty}^{\infty} \delta(t) dt = 1 \quad (1.69)$$

If the $g(t)$ is convoluted with a function $v(t)$ the response of the system for $v(t)$ as input is obtained:

$$v(t) * g(t) = \int_{\tau=0}^{\infty} v(\tau) g(t - \tau) d\tau = f(v(t)) \quad (1.70)$$

[15] By using the convolution theorem:

$$\mathcal{L}\{g(t)\} \mathcal{L}\{v(t)\} = \mathcal{L}\{g * v\}(t) \quad (1.71)$$

[16] the transfer function is defined as:

$$\mathcal{L}\{f(\delta(t))\} \quad (1.72)$$

This definition is usefull since the response of a system for a given input $u(t)$ can be deterimend by the result of the multiplication $G(s)U(s)$. [15] Alternatively if a system is given in state space representation the transfer function can be expressed as:

$$G(s) = C(sE - A)^{-1}B + D \quad (1.73)$$

[14]

1.4.5. Controallability and Observability

For LTI systems in state space representation it can be determend which states are to what degree controllable and observable. This can be computed by the controllability and observability gramians:

$$W_c = \lim_{t \rightarrow \infty} \int_0^t e^{A\tau} B B^* e^{A^*\tau} d\tau \quad (1.74)$$

$$W_o = \lim_{t \rightarrow \infty} \int_0^t e^{A^*\tau} C^* C e^{A\tau} d\tau \quad (1.75)$$

The degree of controlability for a state x can be determined by $x^* W_c x$. If the result of this calculation is large, the system is controlable in the x direction. By swapping W_c with W_o the degree of observability in state x can be computed: $x^* W_o x$. Again if the resulting value is large, the system can be obseverd well in state x . [17]

2. Model Order Reduction

2.1. Introduction

Model Order Reduction (MOR) is a technique to reduce the computational effort of simulate a system using mathematical models. This is done by using two different approaches. The first one is using a numerical approach and the second approach is based on system theory. The first method discussed is called Proper Orthogonal Decomposition. After that Balanced and Modal Truncation will be discussed. The last method discussed will be Hankel-norm approximation. [18].

2.2. Proper Orthogonal Decomposition

The proper orthogonal decomposition (POD) is a method for model order reduction. The reduction in computational effort is done by approximating a solution to a PDE using an orthogonal expansion. The basis functions are obtained by decomposing a set of solutions to the PDE using the SVD.

2.2.1. Orthogonal Expansion

A function $f : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ can be represented by the series:

$$f(x, t) = \sum_{i=1}^{\infty} a_i(t) \phi_i(x) \quad (2.1)$$

$$x, t \in \mathbb{R} \quad (2.2)$$

Here all $\phi(x)$ are orthogonal basis functions.

$$\langle \phi_i, \phi_j \rangle = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{else} \end{cases} \quad (2.3)$$

By using a finite sum instead of the entire series f can be approximated:

$$\tilde{f}(x, t) = \sum_{i=1}^n a_i(t) \phi_i(x) \quad (2.4)$$

Since all ϕ are known, 2.4 has to be solved for the set $\{a_0, \dots, a_n\}$. In case f is known, the solutions contained in this set can be calculated as:

$$a_i(t) = \frac{\langle \phi_i, f \rangle}{\langle \phi_i, \phi_i \rangle} \quad (2.5)$$

[19]

2.2.2. Proper Orthogonal Decomposition for PDEs

Since a PDEs are equations in terms of partial derivatives, the notation $\mathcal{P}(\partial x)u$ is introduced, which denotes a differential operator in terms of the spacial variables x for a function $u(x, t, p)$. The vector p contains some parameters. This is done to provide a more abstract way to denote PDEs:

$$\frac{\partial u}{\partial t} = \mathcal{P}(\partial x)u \quad (2.6)$$

[20] Solving for u is often difficult to impossible. A method that is often used to solve PDEs is called separation of variables. This separation of variables assumes, that the underlying solution $u(x, t)$ can be expressed by 2.2, to solve for $a_k, 0 \leq k \leq n$. Since it is not practical to compute an infinite series, u only gets approximated by using 2.4 instead:

$$\sum_{i=1}^n \frac{\partial a_i}{\partial t} \phi_i = \sum_{i=1}^n \mathcal{P}(\partial x) \phi_i a_i \quad (2.7)$$

By discretizing the spatial dimension along x 2.7 can be expressed using matrix notation:

$$\Phi = [\phi_0, \dots, \phi_n] \quad (2.8)$$

$$\Phi \frac{d}{dt} a = \mathcal{P}(\partial x) \Phi a \quad (2.9)$$

Since the solution u is unknown 2.5 cannot be computed. However the fact, that all ϕ are orthogonal to each other, can be used to solve for all a_k . It can be done by computing the inner product of 2.7 with all basis functions:

$$\langle \sum_{i=1}^n \frac{\partial a_i}{\partial t} \phi_i, \phi_k \rangle = \langle \sum_{i=1}^n \mathcal{P}(\partial x) \phi_i a_i, \phi_k \rangle \quad \forall 0 \leq k \leq n \quad (2.10)$$

This resembles the Galerkin projection. In matrix notation it can be expressed:

$$\Phi^* \Phi \frac{d}{dt} a = \Phi^* \mathcal{P}(\partial x) \Phi a \quad (2.11)$$

By considering 2.3 the equation 2.10 can be expressed as a system of ODEs which can be solved:

$$I \frac{d}{dt} a = \Phi^* \mathcal{P}(\partial x) \Phi a \quad (2.12)$$

After the vector of coefficients a for each time step has been computed, the solution can be assembled:

$$u(x, t) \approx \Phi(x) a(t) \quad (2.13)$$

[21]

2.2.3. Choosing basis vectors

As discussed in the previous section, a set of basis vectors can be used to generate approximate solutions to PDEs. However it was not discussed how those basis functions are chosen. For POD to work, a so called snapshot matrix X has to be available. This snapshot matrix stores

a set of solutions where x_k is the solution for a PDE at time step $k\Delta t$:

$$X = [x_1, \dots, x_m] \quad (2.14)$$

The solutions can be obtained by conducting an experiment on a physical system that is described by the PDE or by simulating the evolution of that PDE. In this paper the solution is obtained using FEM 1.2. SVD is used to decompose the snapshot matrix X . Since the columns of the snapshot matrix contain spacial information at a given point in time and the basis vectors are supposed to encode the spacial information of a solution u the r most dominant left singular values have to be extracted:

$$\tilde{X} = \tilde{U}\tilde{\Sigma}\tilde{V}^* \quad (2.15)$$

$$\Phi = [u_1, \dots, u_r] \quad (2.16)$$

Here the FEM solution is obtained by discretizing the PDE into a large number(n) of spacial nodes. This results in a high dimensional system of ODEs. Since Φ contains only r column ($r \ll n$) vectors a reduction in order can be achieved. However the benefits of this reduced order model are only relevant after the PDE has been solved once using a high dimensional system of ODEs. [21]

2.3. Balanced Truncation

Balanced truncation is a method for model order reduction. The goal of balanced truncation is to approximate a system using only the most relevant modes of the system. The difference to POD is that the modes are not selected by the variance they capture but by the controllability and observability of the modes. This is done by finding a coordinate transform.

2.3.1. Balancing coordinate transform

To find a reduced order model using balanced truncation a orthonormal coordinate transform is applied: $x = Tz$. This yields a new system:

$$\dot{z} = \hat{A}z + \hat{B}u \quad (2.17)$$

$$y = \hat{C}z + Du \quad (2.18)$$

$$\hat{A} = T^{-1}AT \quad \hat{B} = T^{-1}B \quad \hat{C} = CT \quad (2.19)$$

The gramians of this ROM can be obtained by applying 1.75 and 1.74 to 2.19. This yields $\hat{W}_c = T^{-1}W_cT^{-*}$ and $\hat{W}_o = T^*W_oT$ with $T^{-*} := (T^{-1})^* := (T^*)^{-1}$. A requirement T has to satisfy is that it has to make the observability and controllability gramians of the ROM equal and diagonal.

$$\hat{W}_c = \hat{W}_o = \Delta \quad (2.20)$$

$$\hat{W}_c\hat{W}_o = \Delta^2 \quad (2.21)$$

$$T^{-1}W_cW_oT = \Delta^2 \quad (2.22)$$

$$W_cW_oT = T\Delta^2 \quad (2.23)$$

Since Δ is a diagonal matrix, 2.23 is equal to the eigendecomposition of W_cW_o . Therefore T contains the eigenvectors of W_cW_o . However T needs to be rescaled to make \hat{W}_c and \hat{W}_o equal. Here T_u denotes the unscaled eigenvectors of this eigendecomposition that yields gramians that are not equal to each other:

$$T_u^{-1}W_cT_u^{-*} = \Delta_c \quad (2.24)$$

$$T_u^*W_oT_u = \Delta_o \quad (2.25)$$

Scaling T_u by some diagonal matrix Δ_s results in $\Delta_c = \Delta_o$:

$$\Delta_s = \Delta_c^{\frac{1}{4}} \Delta_o^{-\frac{1}{4}} \quad (2.26)$$

$$T = T_u\Delta_s \quad (2.27)$$

Another important property of this transform is that the new coordinates are hierarchically ordered by observability and controllability. It can be shown by deriving some unit vector ζ

that maximizes the controllability and observability:

$$\zeta = \arg \max \zeta^* W_c W_o \zeta \text{ s.t. } \|\zeta\|_2^2 = 1 \quad (2.28)$$

$$\frac{d}{d\zeta} \zeta^* W_c W_o \zeta - 2\lambda \zeta = 0 \quad (2.29)$$

As shown here [22] the remaining derivative can be solved in the following way:

$$\frac{d}{dx} x^* A x = 2Ax \quad (2.30)$$

This holds if A is symmetric. Here both W_c and W_o share the same set of eigenvectors [22], therefore they commute [22]. This means, that the product of W_c and W_o is also symmetric [23]. Applying 2.30 to 2.29 yields:

$$W_c W_o \zeta = \lambda \zeta \quad (2.31)$$

Since λ is the Lagrange multiplier, it is a scalar. It is clear that ζ is a eigenvector. Since the eigenvalues in Δ contain information about how much each eigenvector gets scaled by multiplying it with $W_c W_o$ the eigenvectors can be ordered by controllability and observability.

2.3.2. Mode Truncation

Since the goal of balanced truncation is to find a ROM of rank r that approximates the original system of rank n with $r \ll n$ it is necessary to truncate the balanced system. This yields the following system:

$$\frac{d\tilde{x}}{dt} = \tilde{A}\tilde{x} + \tilde{B}u \quad (2.32)$$

$$y = \tilde{C}\tilde{x} + \tilde{D}u \quad (2.33)$$

The new state vector \tilde{x} is defined as:

$$\tilde{x} = \begin{bmatrix} z_1 \\ \vdots \\ z_r \end{bmatrix} \quad \zeta = \begin{bmatrix} z_{r+1} \\ \vdots \\ z_n \end{bmatrix} \quad z = \begin{bmatrix} \tilde{x} \\ \zeta \end{bmatrix} \quad (2.34)$$

$$T = \begin{bmatrix} \Psi & T_t \end{bmatrix} \quad T^{-1} = S = \begin{bmatrix} \Phi^* \\ S_t \end{bmatrix} \quad (2.35)$$

By substituting 2.34 and 2.35 into the system in 2.17 and 2.18 the system beecomes:

$$\frac{d}{dt} \begin{bmatrix} \tilde{x} \\ \zeta \end{bmatrix} = \begin{bmatrix} \Phi^* A \Psi & \Phi^* A T_t \\ S_t A \Psi & S_t A T_t \end{bmatrix} \begin{bmatrix} \tilde{x} \\ \zeta \end{bmatrix} + \begin{bmatrix} \Phi^* B \\ S_t B \end{bmatrix} u \quad (2.36)$$

$$y = \begin{bmatrix} C \Psi & C T_t \end{bmatrix} \begin{bmatrix} \tilde{x} \\ \zeta \end{bmatrix} + D u \quad (2.37)$$

However the only relevant part of this system is:

$$\frac{d\tilde{x}}{dt} = \Phi^* A \Psi \tilde{x} + \Phi^* B u \quad (2.38)$$

$$y = C \Psi \tilde{x} + D u \quad (2.39)$$

since this is the only part nescesarry to calculate \tilde{x} . [24]

3. Implementation

The implementation of the discussed methods for solving the heat equation and for model order reduction was done in Matlab. Matlab was chosen as the programming language because it natively features matrix multiplication which finds heavy use in the previously mentioned methods. The second reason for this selection is that there exist ToolBoxes that already implement certain model order reduction methods such as MORLAB [25] or MOR toolbox [26]. The following figure shows the class diagramm of the implementation:

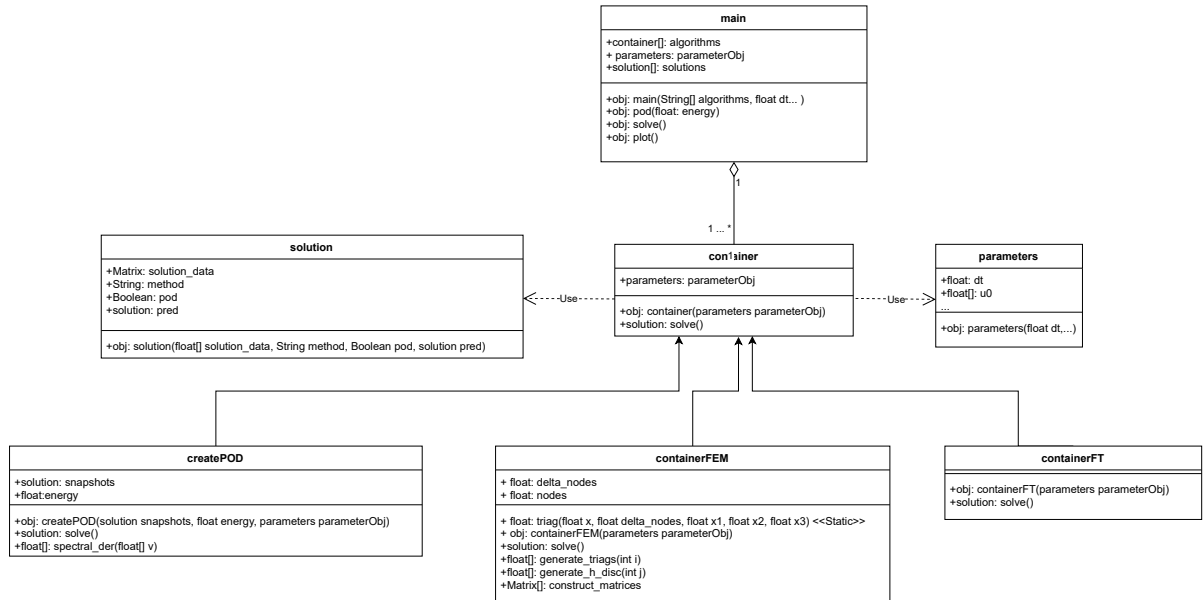


Figure 3.1.: Class diagramm of MOR and FEM implementation

3.1. Class main

The class main is responsible for generating the finite element solution, model order reduction steps and plotting the results. The process can be seen in the following flow chart:

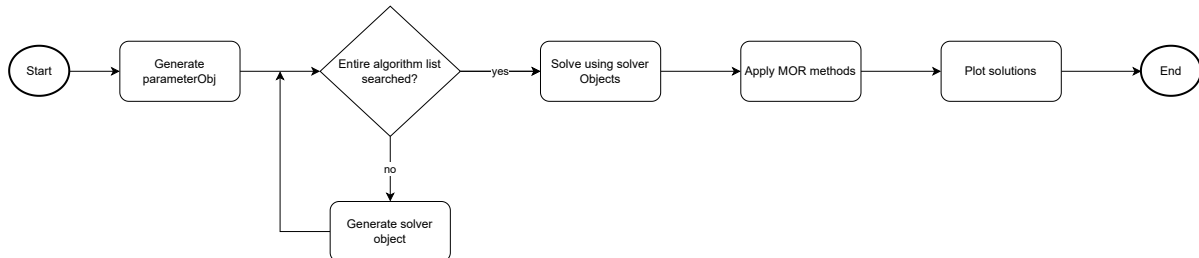


Figure 3.2.: Flow chart of main class

The first step is to generate a parameter object. The parameter object stores all parameters in order to increase the transparency and robustness of the program. The second step is to iterate the array of stated algorithms to solve the heat equation. The options are to solve the heat equation using finite element method 1.2 or using fourier transform 1.1. After all solver objects have been generated, the according solutions are being computed. After that, the MOR methods are displayed. The final step is to display the solutions.

3.2. Class containerFEM

The class containerFEM is responsible for generating a solution using finite element method. FEM is implemented in the following way:

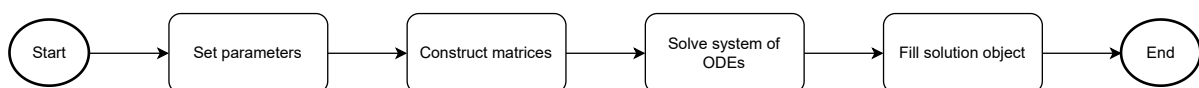


Figure 3.3.: Flow chart for FEM class

The first step is to set the parameters. After that the matrices discussed in 1.2 are being constructed. The next step is to solve the resulting system of ODEs and pass the solution to a solution object.

3.2.1. Construct Matrices

As defined in 1.33 the matrices K and M have to be constructed. Also to compute the initial condition given by u_0 F has to be known 1.35. This is done by the following method:

Algorithm 1 Construct matrices K , M and F

```
1:  $ii \leftarrow \frac{2}{3}\Delta nodes$ 
2:  $ij \leftarrow \frac{1}{6}\Delta nodes$ 
3:  $F \leftarrow \text{zeros}(nodes, 1)$ 
4:  $K \leftarrow \text{zeros}(nodes)$ 
5:  $M \leftarrow \text{zeros}(nodes)$ 
6: for  $i = 1$  to  $nodes$  do
7:    $F(i) \leftarrow \text{trapz}(\phi_i \cdot u_0)$ 
8: end for
9: for  $i = 1$  to  $nodes$  do
10:    $K(i, i) \leftarrow \frac{-2}{\Delta nodes}$ 
11:    $M(i, i) \leftarrow ii$ 
12:   if  $i - 1 > 1$  then
13:      $K(i, i - 1) \leftarrow \frac{1}{\Delta nodes}$ 
14:      $M(i, i - 1) \leftarrow ij$ 
15:   end if
16:   if  $i + 1 < nodes + 1$  then
17:      $K(i, i + 1) \leftarrow \frac{1}{\Delta nodes}$ 
18:      $M(i, i + 1) \leftarrow ij$ 
19:   end if
20: end for
21: return  $[F, K, M]$ 
```

3.2.2. Compute vector $d(t)$

As discussed in 1.2 vector d has to be computed for each time step:

Algorithm 2 Construct vector d

```
1:  $d \leftarrow \text{zeros}(nodes, 1)$ 
2: for  $i = 1$  to  $nodes$  do
3:    $d(i) \leftarrow \text{trapz}(\phi_i \cdot h(x, t_0))$ 
4: end for
5: return  $d$ 
```

The entries of vector d become the integral of the product of a basisfunction and h evaluated at timestep t_0 . This timestep is a argument of that method.

3.2.3. Solve System of ODEs

The most important step in the process of generating a solution is to solve the system of ordinary differential equations that FEM yields:

Algorithm 3 Solve system of ODEs using euler's scheme

```
1:  $[F, K, M] \leftarrow \text{construct\_matrices}()$ 
2:  $C \leftarrow \text{zeros}(\text{nodes}, n\_time\_steps)$ 
3:  $c_0 \leftarrow M^{-1}F$ 
4:  $C(:, 1) \leftarrow c_0$ 
5:  $M(1, :) \leftarrow [0, \dots, 0]$ 
6:  $M(\text{end}, :) \leftarrow [0, \dots, 0]$ 
7:  $N \leftarrow M^{-1}K$ 
8: for  $t = 2$  to  $n\_time\_steps$  do
9:    $d \leftarrow \text{generate\_h\_disc}(t)$ 
10:   $c_n \leftarrow \Delta t N c_0 + M^{-1}h + c_0$ 
11:   $c_0 \leftarrow c_n$ 
12:   $C(:, t) \leftarrow c_n$ 
13: end for
14:  $S \leftarrow []$ 
15: for  $t = 1$  to  $n\_time\_steps$  do
16:   $c \leftarrow C(:, t)$ 
17:   $\text{interpol} \leftarrow \text{interp1}(\text{linspace}(0, L, \text{nodes}), c, X)$ 
18:   $S(:, t) \leftarrow \text{interpol}$ 
19: end for
20: return  $\text{solution}(S, "FEM", 0, 0)$ 
```

In the first line the matrices F , K and M are retrieved. After that the initial vector of coefficients c_0 is computed using LS 1.36 in line three. The next two following lines force the boundary conditions as stated in 1.41. In the lines 8 to 13 1.38 is implemented. In the last step the coefficients are interpolated in spatial direction to fit the given domain X and stored in a solution object.

Literaturverzeichnis

- [1] Gustafsson, B. “Heat Conduction”. In: *Fundamentals of Scientific Computing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 255–262. URL: https://doi.org/10.1007/978-3-642-19495-5_16.
- [2] Papula, L. “Gewöhnliche Differentialgleichungen”. In: *Mathematik für Ingenieure und Naturwissenschaftler Band 2: Ein Lehr- und Arbeitsbuch für das Grundstudium*. Wiesbaden: Springer Fachmedien Wiesbaden, 2015, pp. 343–541. URL: https://doi.org/10.1007/978-3-658-07790-7_4.
- [3] Gustafsson, B. “Finite Element Methods”. In: *Fundamentals of Scientific Computing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 173–192. URL: https://doi.org/10.1007/978-3-642-19495-5_11.
- [4] Gustafsson, B. “Polynomial Expansions”. In: *Fundamentals of Scientific Computing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 117–123. URL: https://doi.org/10.1007/978-3-642-19495-5_7.
- [5] Gustafsson, B. “Least Square Problems”. In: *Fundamentals of Scientific Computing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 125–133. URL: https://doi.org/10.1007/978-3-642-19495-5_8.
- [6] Brunton, S. L./ Kutz, J. N. “Singular Value Decomposition (SVD)”. In: *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. Cambridge University Press, 2019, pp. 3–46.
- [7] Brunton, S. L./ Kutz, J. N. “Fourier and Wavelet Transforms”. In: *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. Cambridge University Press, 2019, pp. 47–83.
- [8] Szabo, F. E. “U”. In: *The Linear Algebra Survival Guide*. Ed. by Szabo, F. E. Boston: Academic Press, 2015, pp. 385–392. URL: <https://www.sciencedirect.com/science/article/pii/B978012409520550028X>.
- [9] Naoko, S. *Geometric Interpretation of the Correlation between Two Variables*. 2018. URL: <https://medium.com/@ns2586/geometric-interpretation-of-the-correlation-between-two-variables-4011fb3ea18e> (visited on 03/01/2023).

- [10] Oliver, K. *Math 19b: Linear Algebra with Probability*. 2011. URL: https://people.math.harvard.edu/~knill/teaching/math19b_2011/handouts/lecture12.pdf (visited on 03/01/2023).
- [11] Victor, L. *IAML: Dimensionality Reduc6on*. 2011. URL: <https://www.inf.ed.ac.uk/teaching/courses/iaml/2011/slides/pca.pdf> (visited on 03/01/2023).
- [12] Douglas, B. “The Fundamentals of Control Theory”. In: Academia, 2019, pp. 8–15.
- [13] Douglas, B. “The Fundamentals of Control Theory”. In: Academia, 2019, pp. 42–47.
- [14] Benner, P. et al. “1 Model order reduction: basic concepts and notation”. In: *Volume 1 System- and Data-Driven Methods and Algorithms*. Ed. by Benner, P. et al. Berlin, Boston: De Gruyter, 2021, pp. 1–14. URL: <https://doi.org/10.1515/9783110498967-001>.
- [15] Douglas, B. “The Fundamentals of Control Theory”. In: Academia, 2019, pp. 48–63.
- [16] Abell, M. L./ Braselton, J. P. “Chapter 8 - Introduction to the Laplace Transform”. In: *Introductory Differential Equations (Fifth Edition)*. Ed. by Abell, M. L./ Braselton, J. P. Fifth Edition. Academic Press, 2018, pp. 399–460. URL: <https://www.sciencedirect.com/science/article/pii/B9780128149485000082>.
- [17] Brunton, S. L./ Kutz, J. N. “CONTROLLABILITY AND OBSERVABILITY”. In: *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. Cambridge University Press, 2019, pp. 335–341.
- [18] *Volume 1 System- and Data-Driven Methods and Algorithms*. Berlin, Boston: De Gruyter, 2021. URL: <https://doi.org/10.1515/9783110498967>.
- [19] Gustafsson, B. “Polynomial Expansions”. In: *Fundamentals of Scientific Computing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 117–123. URL: https://doi.org/10.1007/978-3-642-19495-5_7.
- [20] Gustafsson, B. “Numerical Methods for Differential Equations”. In: *Fundamentals of Scientific Computing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 137–144. URL: https://doi.org/10.1007/978-3-642-19495-5_9.

- [21] Brunton, S. L./ Kutz, J. N. “Reduced Order Models (ROMs)”. In: *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. Cambridge University Press, 2019, pp. 375–402.
- [22] (<https://math.stackexchange.com/users/26325/johannes-kloos>), J. K. *When is matrix multiplication commutative?* Mathematics Stack Exchange. URL:<https://math.stackexchange.com/q/170371>. URL: <https://math.stackexchange.com/q/170371>. (version: 2014-03-13). eprint: <https://math.stackexchange.com/q/170371>.
- [23] Sapir, D. M. V. *Symmetric, diagonal, triangular matrices*. Linear Algebra Web-Notes. <https://math.vanderbilt.edu/sapirmv/msapir/jan22.html#:text=following%20result%20h,Theorem.,BA%20then%20AB%20is%20symmetric>. (version: 2014-03-13). eprint: <https://math.vanderbilt.edu/sapirmv/msapir/jan22.html#:~:text=following%20result%20holds.,Theorem.,BA%20then%20AB%20is%20symmetric..> URL: <https://math.vanderbilt.edu/sapirmv/msapir/jan22.html#:~:text=following%5C%20result%5C%20holds.,Theorem.,BA%5C%20then%5C%20AB%5C%20is%5C%20symmetric..>
- [24] Brunton, S. L./ Kutz, J. N. “BALANCED MODELS FOR CONTROL”. In: *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. Cambridge University Press, 2019, pp. 376–387.
- [25] Benner, P./ Werner, S. W. R. *Morlab*. URL: <https://www.mpi-magdeburg.mpg.de/projects/morlab> (visited on 02/27/2023).
- [26] Systems, P. D. *MOR toolbox*. URL: <https://mordigitalsystems.fr/en/products.html> (visited on 02/27/2023).
- [27] D’Erchie, M. e. a. “Analysis Eins”. In: *Mathe für Nicht-Freaks*. Serlo Education e. V., 2017, pp. 111–113. URL: https://mathe-builds.serlo.org/analysis1_final/analysis1_final.pdf.
- [28] Alexandre M. Bayen, T. S. “An Introduction to MATLAB® Programming and Numerical Methods for Engineers”. In: Science Direct, 2015, pp. 211–214. URL: <https://www.sciencedirect.com/book/9780124202283/an-introduction-to-matlab-programming-and-numerical-methods-for-engineers>.

A. Appendix

A.1. Deriving matrices for FEM using piecewise linear functions

The so called triangle function is defined as follows:

$$\phi_j(x) = \begin{cases} (x - x_{j-1})/\Delta x, & x_{j-1} \leq x < x_j \\ (x_{j+1} - x)/\Delta x, & x_j \leq x < x_{j+1} \\ 0, & \text{otherwise} \end{cases} \quad (\text{A.1})$$

[4] The following integrals have to be evaluated:

$$\int_{\chi} \phi_j \phi_k dx \quad \forall \phi_k \in \phi \quad (\text{A.2})$$

$$- \int_{\chi} \frac{d\phi_j}{dx} \frac{d\phi_k}{dx} dx \quad \forall \phi_k \in \phi \quad (\text{A.3})$$

With $\chi \subset \mathbb{R}$. Note that the product of two functions ϕ_j and ϕ_k and their derivatives is only under two conditions not zero:

1. $k = j$

Considering this case the integral A.2 becomes:

$$\int_{x_{j-1}}^{x_j} \phi_j^2 dx + \int_{x_j}^{x_{j+1}} \phi_j^2 dx \quad (\text{A.4})$$

Because of symmetry only one of the above integrals have to computed:

$$2 \int_{x_{j-1}}^{x_j} \phi_j^2 dx \quad (\text{A.5})$$

$$= \frac{2}{\Delta x^2} \int_{x_{j-1}}^{x_j} (x - x_{j-1})^2 dx \quad (\text{A.6})$$

$$\frac{2}{3\Delta x^2} [(x - x_{j-1})^3]_{x_{j-1}}^{x_j} = \frac{2}{3\Delta x^2} \Delta x^3 = \frac{2}{3} \Delta x \quad (\text{A.7})$$

Integral A.3 for $i = j$ taking symmetry into account becomes:

$$- \int_{x_{j-1}}^{x_{j+1}} \left(\frac{d\phi_j}{dx} \right)^2 dx = - \frac{1}{\Delta x^2} \int_{x_{j-1}}^{x_{j+1}} 1 dx \quad (\text{A.8})$$

$$= - \frac{1}{\Delta x^2} [x]_{x_{j-1}}^{x_{j+1}} = - \frac{2}{\Delta x} \quad (\text{A.9})$$

2. $|j - k| = 1$ A.2 becomes:

$$\frac{1}{\Delta x^2} \int_{x_j}^{x_{j+1}} (x - x_j)(x_{j+1} - x) dx \quad (\text{A.10})$$

$$= \left[\frac{1}{2} x^2 x_{j+1} - \frac{1}{3} x^3 - x x_{j+1} x_j + \frac{1}{2} x^2 x_j \right]_{x_j}^{x_{j+1}} = \frac{1}{6 \Delta x^2} \Delta x^3 = \frac{1}{6} \Delta x \quad (\text{A.11})$$

Finally A.3 has to be evaluated for this condition:

$$- \int_{x_j}^{x_{j+1}} \frac{d\phi_j}{dx} \frac{d\phi_{j+1}}{dx} dx = \frac{1}{\Delta x^2} \int_{x_j}^{x_{j+1}} 1 dx = \frac{1}{\Delta x^2} [x]_{x_j}^{x_{j+1}} = \frac{1}{\Delta x} \quad (\text{A.12})$$

A.2. Proof that matrix M is invertible

Let M_n be a matrix with $M_n \in \mathbb{R}^{n \times n}$ given by:

$$m_{ij} = \begin{cases} a, & k = j \\ b, & |k - j| = 1 \\ 0, & otherwise \end{cases} \quad (\text{A.13})$$

It's determinant is given by the Laplace expansion:

$$\det(M_n) = \sum_{j=1}^n (-1)^{i+j} a_{ij} N_{ij} \quad \forall i \quad (\text{A.14})$$

N_{ij} is the determinant of the matrix M' that is obtained by removing the i^{th} row and j^{th} column of M_n . This expression can be simplified using the definition of m_{ij} :

$$\det(M_n) = aN_{11} - bN_{12} \quad (\text{A.15})$$

N_{11} is equivalent to $\det(M_{n-1})$, since the indices of rows and columns of M' are in consecutive order and M' is a $n-1 \times n-1$ matrix :

$$N_{11} = \det(M') \quad (\text{A.16})$$

$$M' = \begin{bmatrix} m_{22} & \dots & m_{2n} \\ \vdots & \ddots & \vdots \\ m_{n2} & \dots & m_{nn} \end{bmatrix} \quad (\text{A.17})$$

N_{12} can be obtained by calculating the determinant of M' using the Laplace expansion:

$$M' = \begin{bmatrix} m_{21} & m_{23} & \dots & m_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ m_{n1} & m_{n3} & \dots & m_{nn} \end{bmatrix} \quad (\text{A.18})$$

$$\det(M') = m_{21} \cdot \det(M'') \quad (\text{A.19})$$

$$M'' = \begin{bmatrix} m_{33} & \dots & m_{3n} \\ \vdots & \ddots & \vdots \\ m_{n3} & \dots & m_{nn} \end{bmatrix} \quad (\text{A.20})$$

In A.19 only the stated term has to be evaluated since all entries of the first column of the second submatrix are zero. Therefore the determinant is zero. The row and column indices of M'' are in consecutive order and it is a $(n-2) \times (n-2)$ matrix. Therefore M'' is equivalent to M_{n-2} . A.15 becomes:

$$\det(M_n) = a \cdot \det(M_{n-1}) - b^2 \cdot \det(M_{n-2}) \quad (\text{A.21})$$

Furthermore this implies $\det(M_0) = 1$:

$$\det(M_2) = a^2 - b^2 = a \cdot \det(M_1) - b^2 \cdot 1 \quad (\text{A.22})$$

$$\Rightarrow \det(M_0) = 1 \quad (\text{A.23})$$

Using the definition of 1.33 and 1.2.1 this can be seen as the following sequence:

$$a_0 = 1, \quad a_1 = \frac{2\Delta x}{3} \quad (\text{A.24})$$

$$a_{n+1} = \frac{2\Delta x}{3} \cdot a_n - \frac{\Delta x^2}{36} \cdot a_{n-1} \quad (\text{A.25})$$

As described here [27] a recursive sequence converges if it is monotone and has a limit. A proof by induction shows that this sequence is monotone for $n \geq 1$.

Base case:

$$a_2 = \left(\frac{2\Delta x}{3}\right)^2 - \frac{\Delta x^2}{36} = \Delta x^2 \left(\frac{4}{9} - \frac{1}{36}\right) < \frac{2\Delta x}{3} = a_1 \quad (\text{A.26})$$

Induction step: Assuming that $a_k < a_{k-1}$ holds, $a_{k+1} < a_k$ also holds:

$$a_{k+1} = \frac{2\Delta x}{3} \cdot a_k - \frac{\Delta x^2}{36} \cdot a_{k-1} < \frac{2\Delta x}{3} \cdot a_{k-1} - \frac{\Delta x^2}{36} \cdot a_{k-2} = a_k \quad (\text{A.27})$$

The limit of this sequence is as follows:

$$\alpha = \lim_{n \rightarrow \infty} a_{n+1} = \lim_{n \rightarrow \infty} \Delta x \cdot \frac{2}{3} \cdot \lim_{n \rightarrow \infty} a_n - \lim_{n \rightarrow \infty} \Delta x^2 \cdot \frac{1}{36} \cdot \lim_{n \rightarrow \infty} a_{n-1} = 0 \cdot \alpha - 0 \cdot \alpha = 0 \quad (\text{A.28})$$

Since this series is monotone and converges to zero as n goes to infinity, there is no $n \in \mathbb{N}$ for which $a_n = 0$. Therefore the determinant of the matrix M defined in 1.33 is not zero and M is invertible.

A.3. Equivalence of picewise linear polynomials and linear interpolation

A picewise linear polynomial in the form of:

$$u(x, t) = \sum_{j=1}^N c_j(t) \phi_j(x) \quad (\text{A.29})$$

With ϕ_j being defined as A.1 and $c_j : \mathbb{R} \rightarrow \mathbb{R}$ is equivalent to linear interpolation with respect to x :

$$\hat{u}(x, t) = u_j + \frac{(u_{j+1} - u_j)(x - x_j)}{x_{j+1} - x_j} \quad (\text{A.30})$$

for $x_j < x < x_{j+1}$ [28]. This can be shown by evaluating $u(x, t)$ between two neighbouring ϕ and $x_j < x < x_{j+1}$:

$$u(x, t) = \phi_j(x) c_j(t) + \phi_{j+1}(x) c_{j+1}(t) \quad (\text{A.31})$$

$$= \frac{x_{j+1} - x}{\Delta x} c_j(t) + \frac{x - x_j}{\Delta x} c_{j+1}(t) \quad (\text{A.32})$$

$$= \frac{x_{j+1} - x}{\Delta x} u_j + \frac{x - x_j}{\Delta x} u_{j+1} \quad (\text{A.33})$$

$$= \frac{(x_{j+1} u_j - x u_j) + (x u_{j+1} - x_j u_{j+1})}{\Delta x} \quad (\text{A.34})$$

$$= \frac{(u_{j+1} - u_j)x + x_{j+1} u_j - x_j u_{j+1}}{\Delta x} \quad (\text{A.35})$$

$$= \frac{(u_{j+1} - u_j)x + x_j u_j + \Delta x u_j - x_j u_{j+1}}{\Delta x} \quad (\text{A.36})$$

$$= u_j + \frac{(u_{j+1} - u_j)x - x_j(u_{j+1} - u_j)}{\Delta x} \quad (\text{A.37})$$

$$= u_j + \frac{(u_{j+1} - u_j)(x - x_j)}{x_{j+1} - x_j} \quad (\text{A.38})$$